

Average-case Optimal Approximate Circular String Matching

Carl Barton¹, Costas S. Iliopoulos^{1,2}, and Solon P. Pissis^{1*}

¹ Department of Informatics, King's College London, The Strand, London, UK
`{carl.barton,costas.iliopoulos,solon.pissis}@kcl.ac.uk`

² Department of Mathematics & Statistics, University of Western Australia, 35
 Stirling Highway, Perth, Australia

Abstract. Approximate string matching is the problem of finding all factors of a text t of length n that are at a distance at most k from a pattern x of length m . Approximate circular string matching is the problem of finding all factors of t that are at a distance at most k from x or from any of its rotations. In this article, we present a new algorithm for approximate circular string matching under the edit distance model with optimal average-case search time $\mathcal{O}(n(k + \log m)/m)$. Optimal average-case search time can also be achieved by the algorithms for multiple approximate string matching (Fredriksson and Navarro, 2004) using x and its rotations as the set of multiple patterns. Here we reduce the preprocessing time and space requirements compared to that approach.

Keywords: algorithms on automata and words, average-case complexity, average-case optimal, approximate string matching

1 Introduction

In order to provide an overview of our results and algorithms, we begin with a few definitions, generally following [4]. We think of a *string* x of *length* n as an array $x[0 \dots n-1]$, where every $x[i]$, $0 \leq i < n$, is a *letter* drawn from some fixed *alphabet* Σ of size $\sigma = \mathcal{O}(1)$. By a q -*gram* we refer to any string $x \in \Sigma^q$. The *empty string* of length 0 is denoted by ε . A string x is a *factor* of a string y if there exist two strings u and v , such that $y = uxv$. Consider the strings x, y, u , and v , such that $y = uxv$. If $u = \varepsilon$, then x is a *prefix* of y . If $v = \varepsilon$, then x is a *suffix* of y . Let x be a non-empty string of length n and y be a string. We say that there exists an *occurrence* of x in y , or, more simply, that x *occurs in* y , when x is a factor of y . Every occurrence of x can be characterised by a position in y . Thus we say that x occurs at the *starting position* i in y when $y[i \dots i+n-1] = x$. Given a string x of length m and a string y of length $n \geq m$, the *edit distance*, denoted by $\delta_E(x, y)$, is defined as the minimum total cost of operations required to transform one string into the other. For simplicity, we only count the number of edit operations, considering the cost of each to be 1 [15]. The allowed edit operations are as follows:

* Supported by a London Mathematical Society grant (no. 51303).

- *Insertion*: insert a letter in y , not present in x ; (ε, b) , $b \neq \varepsilon$
- *Deletion*: delete a letter in y , present in x ; (a, ε) , $a \neq \varepsilon$
- *Substitution*: replace a letter in y with a letter in x ; (a, b) , $a \neq b$, and $a, b \neq \varepsilon$.

We write $x \equiv_k^E y$ if the edit distance between x and y is at most k . Equivalently, if $x \equiv_k^E y$, we say that x and y have at most k *differences*. We refer to the *standard dynamic programming matrix* of x and y as the matrix defined by $D[i, 0] = i$, $0 \leq i \leq m$, $D[0, j] = j$, $0 \leq j \leq n$

$$D[i, j] = \min \begin{cases} D[i-1, j-1] + (1 \text{ if } x[i-1] \neq y[j-1]) \\ D[i-1, j] + 1 \\ D[i, j-1] + 1 \end{cases}, 1 \leq i \leq m, 1 \leq j \leq n.$$

Similarly we refer to the *standard dynamic programming algorithm* as the algorithm to compute the edit distance between x and y through the above recurrence in time $\mathcal{O}(mn)$. Given a non-negative integer threshold k for the edit distance, this can be computed in time $\mathcal{O}(mk)$ [17]. We say that there exists an *occurrence* of x in y with at most k differences, or, more simply, that x *occurs in* y with at most k differences, when $u \equiv_k^E x$ and u is a factor of y .

A circular string of length n can be viewed as a traditional linear string which has the left- and right-most symbols wrapped around and stuck together in some way. Under this notion, the same circular string can be seen as n different linear strings, which would all be considered equivalent. Given a string x of length n , we denote by $x^i = x[i..n-1]x[0..i-1]$, $0 < i < n$, the i -th *rotation* of x and $x^0 = x$. Consider, for instance, the string $x = x^0 = \text{abababbc}$; this string has the following rotations: $x^1 = \text{bababbca}$, $x^2 = \text{ababbcab}$, $x^3 = \text{babbcaba}$, $x^4 = \text{abbcabab}$, $x^5 = \text{bbcababa}$, $x^6 = \text{bcababab}$, $x^7 = \text{cabababb}$.

This type of structure occurs in the DNA of viruses, bacteria, eukaryotic cells, and archaea. In [9], it was noted that, due to this, algorithms on circular strings may be important in the analysis of organisms with such structure. For instance, circular strings have been studied before in the context of sequence alignment. In [14,5], algorithms for multiple circular sequence alignment were presented. Here we consider the problem of finding occurrences of a pattern x of length m with circular structure in a text t of length n with linear structure. This is the problem of *circular string matching*.

The problem of exact circular string matching has been considered in [16], where an $\mathcal{O}(n)$ -time algorithm was presented. The approach presented in [16] consists of preprocessing x by constructing a *suffix automaton* of the string xx , by noting that every rotation of x is a factor of xx . Then, by feeding t into the automaton, the lengths of the longest factors of xx occurring in t can be found by the links followed in the automaton in time $\mathcal{O}(n)$. In [6], an average-case optimal algorithm for exact circular string matching was presented and it was also shown that the average-case lower bound for single string matching of $\Omega(n \log_\sigma m/m)$ also holds for circular string matching. Very recently, in [3], the authors presented two fast average-case algorithms based on word-level parallelism. The first algorithm requires average-case time $\mathcal{O}(n \log_\sigma m/w)$, where w is the number of bits in the computer word. The second one is based on a

mixture of word-level parallelism and q -grams. The authors showed that with the addition of q -grams, and by setting $q = \Theta(\log_\sigma m)$, an average-case optimal time of $\mathcal{O}(n \log_\sigma m/m)$ is achieved. Indexing circular patterns [12] based on the construction of *suffix tree*—have also been considered.

The aforementioned algorithms for the exact case have the disadvantage that they cannot be applied in a biological context since single nucleotide polymorphisms and errors introduced by wet-lab sequencing platforms might have occurred in the sequences; also it is not clear whether they could easily be adapted to deal with the approximate case. For the rest of the article, we assume that each position in the text t is uniformly randomly drawn from Σ , and consider the following problem.

APPROXIMATECIRCULARSTRINGMATCHING

Input: a pattern x of length m , a text t of length $n > m$, and an integer threshold $k < m$

Output: all factors u of t such that $u \equiv_k^E x^i$, $0 \leq i < m$

Similar to the exact case [6], it can be shown that the average-case lower bound for single approximate string matching of $\Omega(n(k + \log_\sigma m)/m)$ [2] also holds for approximate circular string matching under the edit distance model. Recently, we have presented average-case $\mathcal{O}(n)$ -time algorithms for approximate circular string matching which are also very efficient in practice [1]. In [10], an algorithm with $\mathcal{O}(\frac{nk \log m}{m})$ average-case search time was presented. To achieve average-case optimality, one could use the algorithms for multiple approximate string matching, presented in [8], for matching the $r = m$ rotations of x with $\mathcal{O}(n(k + \log_\sigma rm)/m)$ average-case search time, only if $k/m < 1/2 - \mathcal{O}(1/\sqrt{\sigma})$ and $r = \mathcal{O}(\min(n^{1/3}/m^2, \sigma^{o(m)}))$. Therefore the focus of this article is on a more *direct* algorithm which also improves on the preprocessing time and space complexity.

Our Contribution. In this article, we present a new average-case optimal algorithm for approximate circular string matching, under the edit distance model, that reduces the preprocessing time and space requirements compared to previous algorithms with optimal average-case search time. These savings are around $\mathcal{O}(m^2)$ or more in all cases.

2 Algorithm

In this section, we present our algorithm for approximate circular string matching under the edit distance model. The presented algorithm consists of two distinct schemes: the *searching* scheme, which determines if the currently considered text window potentially has a valid occurrence; in case the window *may* contain a valid occurrence, we are required to check the window for valid occurrences of the pattern or any of its rotations; this is done through the *verification* scheme.

Intuitively, the algorithm considers a *sliding window* of length $m - k$ of the text, and reads q -grams backwards from the end of the window until it is likely to have found enough *differences* to skip the entire window. That is, we wish to

make the probability of a verification being triggered sufficiently unlikely whilst also ensuring we can shift the window a reasonable amount.

The rest of this section is structured as follows. We first present an efficient incremental string comparison technique which forms the basis of the verification scheme. We then present the searching scheme of our algorithm which requires a preprocessing step. In fact, this preprocessing step is similar to the verification scheme. Finally, we show how plugging these schemes together results in a new average-case optimal algorithm for approximate circular string matching.

2.1 Verification scheme

The verification scheme of our algorithm is based on incremental string comparison techniques. First we give an introduction to these techniques; and then explain how we use them in the verification scheme. The incremental string comparison problem was introduced in the pioneering work of Landau *et al* [13]. The authors considered the following problem: given the edit distance between two strings A and B , how can the edit distance between A and $\mathbf{b}B$ or $B\mathbf{b}$ be efficiently derived, where \mathbf{b} is an additional letter. Given a threshold on the number of differences k , they solve this problem and allow prepending and appending of letters in time $\mathcal{O}(k)$ per operation. Later the authors of [11] considered a generalisation of this problem with the aim of computing all maximal gapped palindromes in a string. The problem considered is a generalisation of the incremental string comparison problem considered in [13] as it considers how to efficiently derive the edit distance when prefixes are deleted and letters are prepended to A or B . The solution proposed in [11] also has a time complexity of $\mathcal{O}(k)$ per operation. The solution for the generalised incremental string comparison problem forms the basis of our verification step. The technique lends itself more naturally to circular string matching due to the increased flexibility it provides. We begin by recalling some of the main results from [11] required for our algorithm.

The main idea in both [13] and [11] is the efficient computation of the so-called *h-waves*. In the standard dynamic programming matrix for two strings x and y , we say that a cell $D[i, j]$ is on the diagonal d iff $j - i = d$. For each diagonal, we may have a lowest cell with value h ; if $D[i, j] = h$ and $D[i + 1, j + 1] = h + 1$ then $D[i, j]$ is this cell for diagonal $j - i$. The *h-wave*, for all $0 \leq h \leq k$, is the position of all these cells across all diagonals, that is, a list H_h of length $\mathcal{O}(k)$, where each entry is a pair (i, j) such that $D[i, j] = h$ and $D[i + 1, j + 1] = h + 1$. Note that the i -th wave can only contain entries on diagonal zero and the i diagonals either side of it, so for $0 \leq i \leq k$ every wave has size $\mathcal{O}(k)$. Both incremental string comparison techniques show some bounds on the possible values of the cells on *h-waves* and how to efficiently compute them. These *h-waves* define the entire dynamic programming matrix due to the monotonicity properties of the matrix. For any diagonal d , if we know the position of the lowest cell on d with value h and $h + 1$, then we also know the value of every cell between these two cells: it must be $h + 1$. So given the *h-waves* of the matrix, for all $0 \leq h \leq k$, we have all the information that is in the standard dynamic programming matrix. The key result from our perspective is the following.

Let $\text{cat}(u', u)$ denote the string obtained by concatenating u' and u , where $u, u' \in \Sigma^+$. Let $\text{del}(\alpha, u)$ denote the string obtained by deleting the prefix of length α of u . Let D' denote the standard dynamic programming matrix for strings $\text{cat}(A', A)$ and $\text{del}(t_2, B)$, where $|A'| = t_1$.

Theorem 1 ([11]). *The 0-wave, 1-wave, \dots , and k -wave of matrix D' can be computed in time $\mathcal{O}((t_1 + t_2)k)$.*

If a window of the text triggers a verification then we have a window of length $m - k$ such that there exist some q -grams of the window that occur in x or its rotations with at most k differences in total. When we verify a window, we check for occurrences of pattern x starting at every position in the window. For each position, we may have a factor of length at most $m + k$ representing an occurrence, meaning we must consider a factor w of the text of length $2m$ which we refer to as a *block*. This ensures we avoid missing any occurrences at the $m - k$ starting positions as $(m - k) + (m + k) = 2m$.

For each possible starting position i , $0 \leq i < m - k$, we compute the 0-wave, 1-wave, \dots , and k -wave for x and $w' = w[i \dots 2m - 1]$, the suffix of w starting at position i . To check if we have an occurrence, we must check the k -wave H_k . We iterate through each entry in the k -wave H_k ; and if H_k has missing entries or contains entries on the last row of the matrix, then x occurs in w with at most k differences.

Similarly we can check for the occurrences of the rotations of x using the incremental string comparison techniques. We are now ready to outline the verification scheme, denoted by function **VER**. Given the pattern x of length m , an integer threshold $k < m$, and a block w of length $2m$ of the text t , function **VER** finds all factors u of w such that $u \equiv_k^E x^i$, $0 \leq i < m$. If any diagonal has no entry on the k -wave then that diagonal reached the last row of the matrix with less than k differences; this means x occurs in w with less than k differences.

Function $\text{VER}(x, m, k, w, 2m)$

Compute the edit distance between x and $w' = w[0 \dots 2m - 1]$ with at most k differences using the standard dynamic programming algorithm;

Check for any occurrences using D , and if found, **report** an occurrence at position 0;

foreach $i \in \{1, m - k - 1\}$ **do**

foreach $j \in \{1, m\}$ **do**

 Construct rotation x^j of x by removing the first letter of x^{j-1} and appending it to the end of x^{j-1} ;

 Compute the edit distance between x^j and $w' = w[i \dots 2m - 1]$ using the incremental string comparison techniques;

 Check for any occurrences using H_k , and if found, **report** an occurrence at the current position i being checked;

Lemma 2. *Given the pattern x of length m , an integer threshold $k < m$, and string w of length $2m$, function VER requires time $\mathcal{O}(m^2k)$.*

Proof. Computing the edit distance between x and $w[0..2m-1]$ with at most k differences takes time $\mathcal{O}(mk)$ using the standard dynamic programming algorithm. By Theorem 1, computing the edit distance between all the rotations of the pattern and $w[i..2m-1]$ for a single position in w requires $\mathcal{O}(mk)$; and there are $\mathcal{O}(m)$ positions in w . In total, the time is $\mathcal{O}(mk + m^2k)$, that is $\mathcal{O}(m^2k)$. \square

2.2 Searching scheme

The searching scheme of the presented algorithm requires the preprocessing and indexing of the pattern x . We first present the preprocessing required and then present the searching technique itself.

Preprocessing. We build a q -gram index in a similar way as that proposed by Chang and Marr in [2]. Intuitively, we wish to determine the minimum possible edit distance between every q -gram and any factor of x or its rotations. Equivalently we find the minimum possible edit distance between every q -gram and any *prefix* of a factor of length $2q$ of x and the suffixes of length 1 to $2q$ of x or its rotations. An index like this allows us to lower bound the edit distance between a window of the text and x or its rotations without computing the edit distance between them. To build this index, we generate every string of length q on Σ , and find the minimum edit distance between it and all prefixes of factors of length $2q$ of x or its rotations. This information can easily be stored by generating a numerical representation of the q -gram and storing the minimum edit distance in an array at this location. If we know the numerical representation, we can then look up any entry in constant time.

We determine the edit distance using the preprocessing scheme, denoted by function PRE , which is similar to the verification scheme (function VER).

Given the string $x' = x[0..m-1]x[0..m-2]$ of length $2m-1$, function PRE finds the minimum edit distance between every q -gram on Σ , generated in increasing order, and any factor u of length $2q$ of x' and its suffixes of length 1 to $2q$.

Lemma 3. *Given the string $x' = x[0..m-1]x[0..m-2]$ of length $2m-1$ on Σ , $\sigma = |\Sigma|$, and $q < m$, function PRE requires time $\mathcal{O}(\sigma^qmq)$ and space $\mathcal{O}(\sigma^q)$.*

Proof. The time required for initialising array M is $\mathcal{O}(\sigma^q)$. The time required for computing the edit distance between $x'[0..2q-1]$ and s is $\mathcal{O}(q^2)$ using the standard dynamic programming algorithm. By Theorem 1, computing the edit distance between all $2q$ -grams of x' and s requires time $\mathcal{O}(mq)$. There exist $\mathcal{O}(\sigma^q)$ possible q -grams on Σ and so, in total, the time complexity is $\mathcal{O}(\sigma^qmq)$. Keeping array M in memory requires space $\mathcal{O}(\sigma^q)$. \square

```

Function  $PRE(x', 2m - 1, q, \sigma)$ 
   $M[0 \dots \sigma^q - 1] \leftarrow 0;$ 
   $j \leftarrow 0;$ 
  foreach  $s \in \Sigma^q$  do
    Compute the edit distance between  $u = x'[0 \dots 2q - 1]$  and  $s$  using
    the standard dynamic programming algorithm. Set  $E_{\min}$  equal to
    the minimum edit distance between  $s$  and any prefix of  $u$  using  $D$ ;
    foreach  $i \in \{1, 2m - q - 1\}$  do
       $u \leftarrow x'[i \dots \min\{i + 2q - 1, 2m - 2\}];$ 
      Compute the edit distance  $E'$  between  $u$  and  $s$  using the
      incremental string comparison techniques. Set  $E'$  equal to the
      minimum edit distance between  $s$  and any prefix of  $u$  using  $H_q$ ;
      if  $E' < E_{\min}$  then  $E_{\min} \leftarrow E'$ 
     $M[j] \leftarrow E_{\min};$ 
     $j \leftarrow j + 1;$ 
  return  $M;$ 

```

Searching. In the search phase we wish to read backwards enough q -grams from a window of size m that the probability we must verify the window is small and the amount we can shift the window by is sufficiently large. We now recall some important lemmas from [2] that we will use in the analysis of our algorithm.

Lemma 4 ([2]). *The probability that two q -grams on Σ , one being uniformly random, have a common subsequence of length $(1 - c)q$ is at most $\frac{a\sigma^{-dq}}{q}$, where $a = (1 + o(1))/(\pi c(1 - c))$ and $d = 1 - c + 2c \log_{\sigma} c + 2(1 - c) \log_{\sigma}(1 - c)$. The probability decreases exponentially for $d > 0$, which holds if $c < 1 - \frac{e}{\sqrt{\sigma}}$.*

Lemma 5 ([2]). *If s is a q -gram occurring with less than cq differences in a given string u , $|u| \geq q$, s has a common subsequence of length $q - cq$ with some q -gram of u .*

By Lemmas 4 and 5, we know that the probability of a random q -gram occurring in a string of length m with less than cq differences is no more than $ma\sigma^{-dq}/q$ as we have $m - q + 1$ q -grams in the string. For circular string matching this is not sufficient. To ensure that we have the q -grams of all possible rotations of pattern x , we instead consider the string $x' = x[0 \dots m - 1]x[0 \dots m - 2]$ and extract the q -grams from x' . We may have up to $2m - q$ q -grams, but to simplify the analysis we assume we have $2m$ and so the probability becomes $2ma\sigma^{-dq}/q$.

In the case when we read $k/(cq)$ q -grams, we know that with probability at most $(k/(cq))2ma\sigma^{-dq}/q$ we have found less than k differences. This does not permit us to discard the window if all q -grams occur with at most cq differences. To fix this, we instead read $1 + k/(cq)$ q -grams. If any q -gram occurs with less than cq differences, we will need to verify the window; but if they all occur

with at least cq differences, we must exceed the threshold k and can shift the window. When shifting the window we have the case that we shift after verifying the window and the case that the differences exceed k so we do not verify the window. If we have verified the window, we can shift past the last position we checked for an occurrence: we can shift by $m - k$ positions. If we have not verified the window, as we read a fixed number of q -grams, we know the minimum-length shift we can make is one position past this point. The length of this shift is at least $m - k - (q + k/c)$ positions. This means we will have at most $\frac{n}{m - k - (q + k/c)} = \mathcal{O}(\frac{n}{m})$ windows. The previous statement is only true assuming $m - q > k + k/c$, as then the denominator is positive. From there we see that we also have the condition that $q + k + k/c$ can be at most ϵm , where $\epsilon < 1$, so the denominator will be $\mathcal{O}(m)$. This puts a slightly stricter condition on c , that is, $c > \frac{k}{\epsilon m - q - k}$.

We can see that, for each window, we verify with probability at most $(1 + k/(cq))2ma\sigma^{-dq}/q$, where $a = (1 + o(1))/(2\pi c(1 - c))$ and $d = 1 - c + 2c \log_\sigma c + 2(1 - c) \log_\sigma(1 - c)$. So the probability that a verification is triggered is

$$\frac{(1 + k/(cq))2ma\sigma^{-dq}}{q}.$$

Because by Lemma 2, verification takes time $\mathcal{O}(m^2k)$, then per window, the expected cost is

$$\frac{(1 + k/(cq))2ma\sigma^{-dq}\mathcal{O}(m^2k)}{q} = \mathcal{O}\left(\frac{(q + k)m^3ka\sigma^{-dq}}{q^2}\right).$$

We wish to ensure that the probability of verifying a window is small enough that the average work done is no more than the work we must do if we skip a window without verification. When we do not verify a window, we read $1 + k/(cq)$ q -grams and shift the window. This means that we read $q + k/c = \mathcal{O}(q + k)$ letters. So a sufficient *condition* is the following:

$$\frac{(q + k)m^3ka\sigma^{-dq}}{q^2} = \mathcal{O}(q + k).$$

Or equivalently the below expression, where f is the constant of proportionality:

$$\frac{(q + k)m^3ka\sigma^{-dq}}{q^2} \leq f(q + k).$$

By rearranging and setting $f = \sigma$ we get the condition on the value of q below:

$$q \geq \frac{3 \log_\sigma m + \log_\sigma k + \log_\sigma a - 2 \log_\sigma q}{d}.$$

From the condition on q we can see that it is sufficient to pick $q = \Theta(\log_\sigma km)$, so asymptotically on m we get the following:

$$q \geq \frac{3 \log_\sigma m + \log_\sigma k - \mathcal{O}(\log_\sigma \log_\sigma km)}{d}.$$

Therefore, for sufficiently large m , the below condition is sufficient for optimality, where $d = 1 - c + 2c \log_\sigma c + 2(1 - c) \log_\sigma(1 - c)$:

$$q = \frac{3 \log_\sigma m + \log_\sigma k}{d}.$$

For this analysis to hold we must be able to read the required number of q -grams to ensure the probability of verifying a window is small enough to negate the work of doing it. Note that the above probability is the probability that at least one of q -grams match with less than cq differences. To ensure we have enough unread random q -grams in the window for Lemma 5 to hold in the above analysis the window must be of size $m - k \geq 2q + 2k/c$. Now we consider the case where $2q + 2k/c > m - k \geq 2q + k/c$. If we have just verified a window then we have enough new random q -grams and our analysis holds. If we have just shifted then we know that all the q -grams we previously read matched with at least cq differences and we have between 1 and k/qc q -grams and the probability that one of these matches with less than cq difference is less than in the analysis above so it holds.

The condition $m - k \geq 2q + k/c$ implies a condition on c , it must be the case that $c \geq \frac{k}{m-k-2q}$. This condition on c is weaker than our previous condition on c , so to determine the *error ratio* $\frac{k}{m}$, we use the stronger condition. Additionally, from Lemma 4, we know that $c < 1 - \frac{e}{\sqrt{\sigma}}$. So we must pick a value for c subject to $\frac{k}{\epsilon m - k - q} \leq c < 1 - \frac{e}{\sqrt{\sigma}}$. This inequality implies a limit on the error ratio for which our algorithm is optimal. Clearly it must be the case that $\frac{k}{\epsilon m - k - q} < 1 - \frac{e}{\sqrt{\sigma}}$ for $\epsilon < 1$. Rearranging the inequality implies the following sufficient condition on our error ratio:

$$\frac{2k}{m} < \epsilon - \frac{q}{m} - \frac{\epsilon e}{\sqrt{\sigma}} + \frac{qe}{m\sqrt{\sigma}} + \frac{ke}{m\sqrt{\sigma}}.$$

From here we can factorise and divide everything by 2 to get the following:

$$\frac{k}{m} < \frac{\epsilon}{2} - \frac{q}{2m} - \frac{e}{2\sqrt{\sigma}} \left(\epsilon - \frac{q}{m} - \frac{k}{m} \right).$$

So asymptotically on m we have:

$$\frac{k}{m} < \frac{\epsilon}{2} - \mathcal{O}\left(\frac{1}{\sqrt{\sigma}}\right).$$

Note that this technique can work for any ratio which satisfies $\frac{k}{m} < \frac{1}{2} - \mathcal{O}\left(\frac{1}{\sqrt{\sigma}}\right)$. For any ratio below this, pick a large enough value for ϵ such that asymptotically on m the algorithm will work in the claimed search time. By choosing a suitable value for c and $q \geq \frac{3 \log_\sigma m + \log_\sigma k}{d}$ we obtain the following result.

Theorem 6. *The problem APPROXIMATECIRCULARSTRINGMATCHING can be solved in optimal average-case search time $\mathcal{O}(n(k + \log_\sigma m)/m)$.*

3 Comparison with Existing Algorithms

To the best of our knowledge, the only other algorithms to achieve optimal average-case search time for approximate circular string matching are the algorithms presented in [8] for multiple approximate string matching. In the analysis of the algorithms in [8] it is assumed that all patterns are random. In [7] the authors re-analyse their algorithms for the problem of circular string matching with the same preprocessing and space costs. In this section, we analyse these results and compare them with our own. We refer to the algorithm presented in Section 2 as BIP. Due to the constant c in the value of q from Lemma 4, the exact preprocessing and space costs for these algorithms depend on the chosen value for c . It is however possible to determine the minimum savings we make based on the value of q used in all algorithms.

Applying the algorithms in [8] to approximate circular string matching requires a reduction to multiple approximate string matching for matching the m rotations of x . The first algorithm in [8] has the following time complexity:

$$\mathcal{O}(n(k + \log_{\sigma} rm)/m).$$

By setting $r = m$ this matches our search time and the result is valid when $k/m < 1/2 - \mathcal{O}(1/\sqrt{\sigma})$, $r = \mathcal{O}(\min(n^{1/3}/m^2, \sigma^{o(m)}))$, and we have $\mathcal{O}(\sigma^q)$ space available, where q is subject to the constraint:

$$q \geq \frac{4 \log_{\sigma} m + 2 \log_{\sigma} r}{d}.$$

Again by setting $r = m$ this becomes $q \geq \frac{6 \log_{\sigma} m}{d}$ and the preprocessing time is $\mathcal{O}(\sigma^q m^2)$. We will refer to this algorithm as FN1. The second algorithm, presented in [8], has the same preprocessing cost and requires space $\mathcal{O}(\sigma^q m)$. We will refer to this algorithm as FN2. The important difference between FN1 and FN2 comes in the condition on q which is slightly lower for FN2:

$$q \geq \frac{3 \log_{\sigma} m + \log_{\sigma} r + \log_{\sigma}(m + \log_2 r)}{d}.$$

Again, setting $r = m$ this becomes:

$$q \geq \frac{4 \log_{\sigma} m + \log_{\sigma}(m + \log_2 m)}{d}.$$

To simplify the comparison between these approaches, we will ignore the factor of $\log_2 m$, and simply say that the value of q for algorithm FN2 is greater than or equal to $\frac{5 \log_{\sigma} m}{d}$. This is lower than the sufficient requirement, so any saving we make using this value must be at least as good or better in reality.

First let us consider FN1. The preprocessing requirement of BIP is $\mathcal{O}(\sigma^q m q)$, so before any savings made due to the value of q for BIP, we have reduced the preprocessing cost by a factor of $\mathcal{O}(\frac{m}{q})$. Given the condition on q for BIP, it is clear that even in the worst case, when $k = \mathcal{O}(m)$, BIP will make a saving of at least $2 \log_{\sigma} m$ on the value of q . This corresponds to an additional saving of

$\mathcal{O}(m^2)$ in preprocessing time bringing the total to $\mathcal{O}(\frac{m^3}{q})$ and $\mathcal{O}(m^2)$ in space. In the case of FN2, we make a saving of at least $\log_\sigma m$ on the value of q . This corresponds to a total saving of $\mathcal{O}(\frac{m^2}{q})$ in preprocessing time and $\mathcal{O}(m^2)$ in space. It should be noted that this is a pessimistic analysis of the savings as we have assumed $k = \mathcal{O}(m)$ and $d = 1$, although it must hold that $d < 1$. Note that the standard dynamic programming algorithm can be used with runtime $\mathcal{O}(m^3)$ for verification and $\mathcal{O}(\sigma^q m q^2)$ for preprocessing. The speed-ups mentioned in the previous section remain significant as we assumed that $k = \mathcal{O}(m)$. We still achieve a preprocessing speed up of at least $\mathcal{O}(m^2)$ and $\mathcal{O}(m)$ against FN1 and FN2, respectively. Table 1 corresponds to this analysis.

Table 1. Comparison of average-case optimal approximate circular string matching algorithms

| Algorithm | Error Ratio (k/m) | Space | Preprocessing Time | Condition on q |
|-----------|--|---------------------------|-----------------------------|--|
| FN1 | $\frac{1}{2} - \mathcal{O}(\frac{1}{\sqrt{\sigma}})$ | $\mathcal{O}(\sigma^q)$ | $\mathcal{O}(\sigma^q m^2)$ | $\frac{6 \log_\sigma m}{d}$ |
| FN2 | $\frac{1}{2} - \mathcal{O}(\frac{1}{\sqrt{\sigma}})$ | $\mathcal{O}(\sigma^q m)$ | $\mathcal{O}(\sigma^q m^2)$ | $\frac{4 \log_\sigma m + \log_\sigma (m + \log_2 m)}{d}$ |
| BIP | $\frac{1}{2} - \mathcal{O}(\frac{1}{\sqrt{\sigma}})$ | $\mathcal{O}(\sigma^q)$ | $\mathcal{O}(\sigma^q m q)$ | $\frac{3 \log_\sigma m + \log_\sigma k}{d}$ |

4 Final Remarks

In this article, we presented a new average-case optimal algorithm for approximate circular string matching. To the best of our knowledge, this algorithm is the first average-case optimal algorithm specifically designed for this problem. Other average-case optimal algorithms exist but with higher preprocessing and space requirements than the presented algorithm. Additionally the considered problem is solved in a more direct fashion, that is, with no reduction to multiple approximate string matching by taking greater advantage of the similarity of the rotations of the pattern.

Our immediate target is twofold:

- first, we plan on tackling the problem of multiple approximate circular string matching. We will try to generalise the approach we have taken here to see if it leads to an average-case optimal algorithm in this case.
- second, we plan on implementing the presented algorithm. We will then compare the respective implementation to other average- and worst-case approaches.

References

1. Barton, C., Iliopoulos, C.S., Pissis, S.P.: Fast algorithms for approximate circular string matching. *Algorithms for Molecular Biology* 9(1), 9 (2014), <http://www.almob.org/content/9/1/9>

2. Chang, W.I., Marr, T.G.: Approximate string matching and local similarity. In: Proceedings of the 5th Annual Symposium on Combinatorial Pattern Matching. pp. 259–273. CPM '94, Springer-Verlag, London, UK (1994)
3. Chen, K.H., Huang, G.S., Lee, R.C.T.: Bit-Parallel Algorithms for Exact Circular String Matching. *The Computer Journal* (2013)
4. Crochemore, M., Hancart, C., Lecroq, T.: *Algorithms on Strings*. Cambridge University Press, New York, NY, USA (2007)
5. Fernandes, F., Pereira, L., Freitas, A.T.: CSA: An efficient algorithm to improve circular DNA multiple alignment. *BMC Bioinformatics* 10(1), 1–13 (2009)
6. Fredriksson, K., Grabowski, S.: Average-optimal string matching. *Journal of Discrete Algorithms* 7(4), 579–594 (2009)
7. Fredriksson, K., Mäkinen, V., Navarro, G.: Flexible music retrieval in sublinear time. *International Journal of Foundations of Computer Science* 17(06), 1345–1364 (2006), <http://www.worldscientific.com/doi/abs/10.1142/S0129054106004455>
8. Fredriksson, K., Navarro, G.: Average-optimal single and multiple approximate string matching. *Journal of Experimental Algorithmics* 9 (Dec 2004), <http://doi.acm.org/10.1145/1005813.1041513>
9. Gusfield, D.: *Algorithms on Strings, Trees and Sequences*. Cambridge University Press (1997)
10. Hirvola, T., Tarhio, J.: Approximate online matching of circular strings. In: Gudmundsson, J., Katajainen, J. (eds.) *Experimental Algorithms*, Lecture Notes in Computer Science, vol. 8504, pp. 315–325. Springer International Publishing (2014)
11. Hsu, P.H., Chen, K.Y., Chao, K.M.: Finding all approximate gapped palindromes. In: Dong, Y., Du, D.Z., Ibarra, O. (eds.) *Algorithms and Computation*, Lecture Notes in Computer Science, vol. 5878, pp. 1084–1093. Springer Berlin Heidelberg (2009)
12. Iliopoulos, C.S., Rahman, M.S.: Indexing circular patterns. In: Proceedings of the 2nd International Conference on Algorithms and Computation. pp. 46–57. WALCOM'08, Springer-Verlag, Berlin, Heidelberg (2008)
13. Landau, G.M., Myers, E.W., Schmidt, J.P.: Incremental string comparison. *SIAM Journal of Computing* 27–2, 557–582 (1998)
14. Lee, T., Na, J.C., Park, H., Park, K., Sim, J.S.: Finding optimal alignment and consensus of circular strings. In: Proceedings of the 21st annual Conference on Combinatorial Pattern Matching. pp. 310–322. CPM'10, Springer-Verlag, Berlin, Heidelberg (2010)
15. Levenshtein, V.I.: Binary codes capable of correcting deletions, insertions, and reversals. *Tech. Rep.* 8 (1966)
16. Lothaire, M.: *Applied Combinatorics on Words*. Cambridge University Press (2005)
17. Ukkonen, E.: On approximate string matching. In: Karpinski, M. (ed.) *Foundations of Computation Theory*, Lecture Notes in Computer Science, vol. 158, pp. 487–495. Springer Berlin Heidelberg (1983)